

# Financial Data Analysis with Apache Spark

SIMONE NICOSANTI, 0334319, University of Rome Tor Vergata, Italy

This document presents strategies, system architecture and empirical results of financial data analysis with Apache Spark.

## 1 INTRODUCTION

The aim of the project is to do batch processing on a financial dataset obtained by Fintech Infront Financial Technology, using Apache Spark as data processing framework. This dataset is about trades on three main european markets during the week from november 8th 2021 to november 14th 2021. We considered a reduced version of this dataset: this version consists of 4 million events (compared with 289 million of original dataset) considering 500 equities and indices on european markets of Paris (FR), Amsterdam (NL) e Frankfurt/Xetra (ETR). Events are registered as they have been received and there may be events without payload.

## 2 SYSTEM ARCHITECTURE

As for system architecture we decided to emulate it with Docker Compose on a single node. We have four main part of the system:

- Client. The Client has been written in Python using PySpark library. An other library used in the client is Redis client library for python, in order to write results on Redis after processing.
- Apache Hadoop. Apache Hadoop has been used as DFS for dataset retrieve and result storage in a distributed way. As for dataset upload to HDFS, it is downloaded from web ([dataset](#)) using `curl` command when Namenode container is started, saved in a local file and then uploaded to HDFS. As for results retrieve and storage, they have been done using Apache Spark. We decided to use only one Hadoop datanode, for the purpose of giving more resources to Spark workers.
- Apache Spark. Apache Spark has been used to both pre processing and processing. With regard to pre processing, we used DataFrame API to clean up the dataset from raw text header and not valid lines, while as for processing, we used both RDD and Spark SQL in order to study differences in execution time. Before result saving, RDD have been converted to DataFrame in order to sort them and save in more readable way
- Redis. Redis has been used as key-value storage to store query results. We decided to store only results obtained using RDD API to reduce storage times. All results have been saved using a JSON as value, with the aim of retrieving them all in one with single Redis GET, and query name as key.

After data processing, results have been also saved on a client local directory `/Results` which is a Docker Volume attached to client container with the purpose of having results on the host machine, too.

## 3 PRE PROCESSING

With regard to pre processing we did it with Spark. Unfortunately, due to original dataset format, we could not directly read it as CSV file because there was a raw text header which was misinterpreted by Spark. As a consequence we decided to:

- Read as TextFile, obtaining a DataFrame of a single column where each row was a CSV line

---

Author's address: Simone Nicosanti, 0334319, snicosanti@gmail.com, simone.nicosanti.01@students.uniroma2.eu, University of Rome Tor Vergata, Rome, Italy.

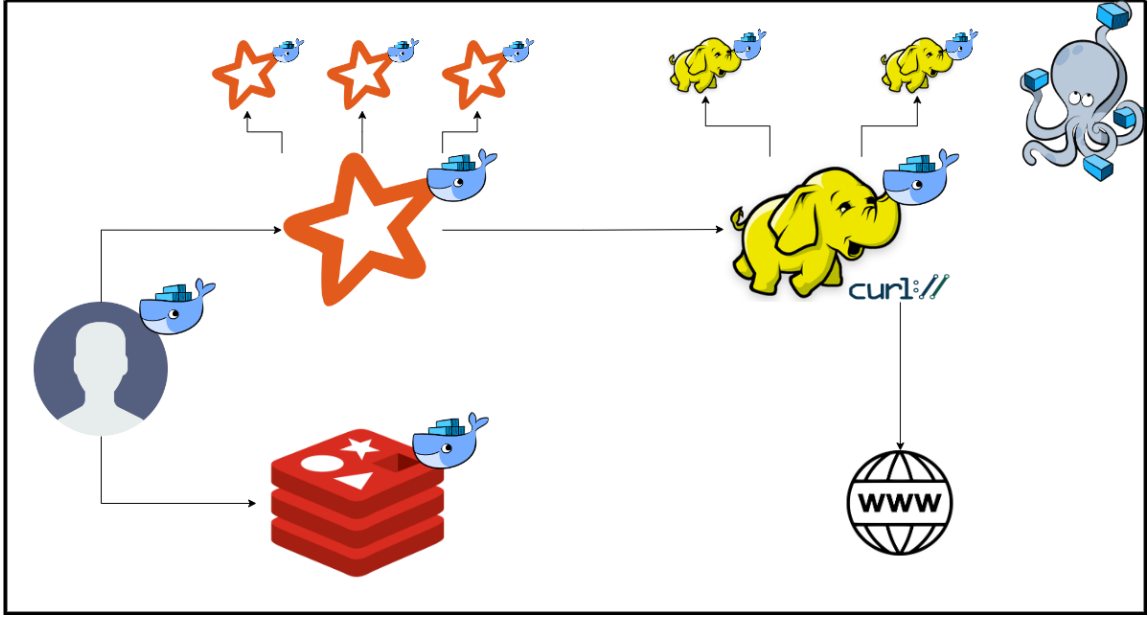


Fig. 1. System Architecture

- Use *split* over the DataFrame splitting by ",", obtaining a DataFrame with all original columns
- Select only interesting columns of CSV

Considering that all queries referenced at most five fields, which are *TradingDate*, *TradingTime*, *ID*, *SecType* and *Last*, we decided to remove all other fields in order to not burden Spark workers. Considering that all queries were based on date and time, we decided to remove all rows which had these fields invalid. After several dataset analysis we noticed that:

- All rows with time equals to *00:00:00.0000* had *Last* field equals to 0. We decided to consider these rows as reset rows for the price and we removed them.
- There were rows which, for the same Date, Time, ID and SecType had different values for the *Last* field: we decided to resume these rows in a single row with *Last* as the average of *Last* fields.

Considering that first and second queries were based based on time hour, we decided to add to DataFrame a redundant column with event hour in order not to compute it again during query processing.

After pre processing operation, we writed dataframe back on HDFS as a *Parquet* file and then read it again as a dataframe before query evaluation. This DataFrame is also converted to RDD and we persisted both in order not to compute them again for query evaluation.

The final schema for both RDD and DataFrame is as follows:

TradingDate	TradingTime	TradingHour	ID	SecType	Last
-------------	-------------	-------------	----	---------	------

Table 1. Final Schema

## 4 QUERY 1

**Query** For each hour, calculate the minimum, average and maximum value of the selling price (Last field) for the only stocks (SecType field with value equal to E) traded on the Paris (FR) markets. In the output also indicate the total number of events used to calculate the statistics. Attention should be paid to events with no payload.

### 4.1 Query 1 with RDD

As for this query using RDD, we did the following steps:

- (1) Filter by ID and SecType
- (2) Map to a Pair RDD with key *(Date, Hour, ID)*
- (3) AggregateByKey using Spark *StatCounter* which returns statistics of values by key obtaining an RDD like *((Date, Hour, ID) (Stats))*
- (4) Map to extract statistics of interest from *StatCounter* result

### 4.2 Query 1 with Spark SQL

As for this query using SQL we did a simple select filtering by *ID* and *SecType*, grouping by *TradingDate*, *TradingHour* and *ID* and then using aggregate operators *min*, *avg*, *max* and *count*.

## 5 QUERY 2

**Query** For each day, for each stock traded on any market, calculate the mean value and standard deviation of their sales price change calculated over a one-hour time window. After calculating the statistical indices on a daily basis, determine the ranking of the best 5 stocks that recorded the best price change during the day and the 5 worst stocks that recorded the worst price change. In the output also indicate the number total number of events used to calculate the statistics.

The hardest part of this query was identify the price of each stock at every hour. After a dataset analysis we found out that there was no stock which had a row at time *hh:00:00.0000* so we decided to consider as the price at time *hh:00* the most recent price of a time before *hh:00*. In addition we did not consider variation equals to 0 for those hour windows when there were no trades in order not to have too many 0 values which would have flattened all statistics: for example if there is a trade at 11:58:43.000 with price of 12.43 and then no other trade until 16:45:32:00, we consider 12.43 as the price of 16:00, in order not to have change equals to zero for all hour windows between 12:00 and 15:00.

### 5.1 Query 2 with RDD

We decided to implement two versions of this Query using RDD, one with a *join* between RDD, the other using a map and a for-cycle.

As for the first variant we have the following steps:

- (1) Map to a key-value RDD of type *((Date, ID, Hour) (TradingTime, Last, TradingTime, Last))*.
- (2) ReduceByKey in order to find first and last price of the stocks in that time window.
- (3) Map to a key-value RDD of type *((Date, ID) (Hour, Last))* obtaining a for each hour window the first and last Trade info.
- (4) Map to a key-value RDD to obtain *partialRDD* which gives the first price for each hour window: this has been achieved summing 1 to *TradingHour*

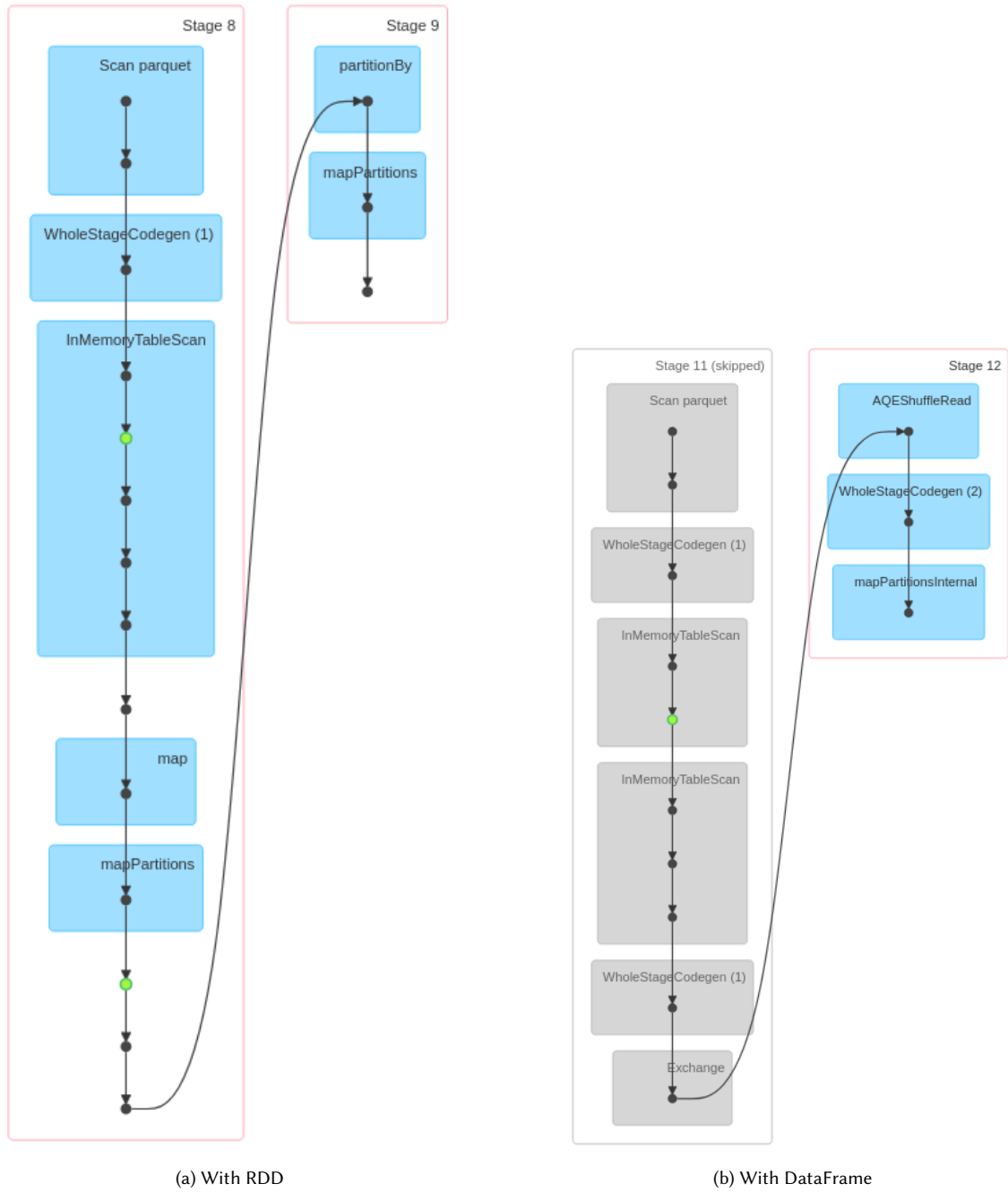


Fig. 2. DAGs for Query 1

- (5) Join *partialRDD* within itself in order to find for a  $(Date, ID)$  couple all possible couples of  $(Time, Last)$ :  $((Date, ID), ((Hour1, Last1) (Hour2, Last2)))$
- (6) Filter for those couples having  $Hour1 < Hour2$  as we are looking for the greatest time before  $Hour2$
- (7) Map to get an RDD  $(Date, ID, Hour2) (Hour1, Last1, Last2)$
- (8) ReduceByKey with key  $(Date, ID, Hour2)$  in order to find the bigger  $Hour1$  before  $Hour2$  and its relative  $Last$ : after this we achieved an RDD of type  $((Date, ID, Hour2), (maxHourBefore, lastBefore, last2))$
- (9) Map to find price variation for each time window obtaining an RDD of  $(Date, ID), variation$
- (10) AggregateByKey in order to find the Variations statistics with Spark StatCounter
- (11) Map in order to extract only interesting statistics from StatCounter result. As we need to find the number of tuples used to compute the statistics we cannot directly use the result of *Count* of StatCounter because it counts the number of values used; nevertheless we can use this number and sum 1 to it: in fact, the number of variations used differs from the number of original tuples only by the first tuple of the list of variations. We now have an RDD of  $(Date, ID), (mean, stdDev, count)$
- (12) Map to an RDD like  $(Date, (mean, stdDev, count, ID))$  as we need to rank by date
- (13) GroupByKey to obtain an RDD like  $Date, listOf((mean, stdDev, count, ID))$
- (14) Maps in order to sort the list for each date and get top 5 and bottom 5. In this way we have a list whose length is of 500 elements at most (which is the max number of equities considered by the dataset); this list is sorted and then we extract some elements. Considering extraction cost as constant we have to consider sorting cost which is  $O(n * \log(n))$ .  
An other way to find the rank would have been looking for max and remove it for five times and similarly for the minimum: this would have led to a cost of  $O(10 * n)$  which is grater compared to  $O(n * \log(n))$  for  $n = 500$
- (15) FlatMap the result in order to obtain again a key-value RDD of type  $((Date, ID) (Avg, StdDev, Count))$

With regard to second variant we have the following:

- (1) First three steps are the same as before
- (2) GroupByKey in order to find an RDD of  $(Date, ID), listOf((Time, Last))$
- (3) Map using a custom function which iterates over the sorted value of each key and computes price variations.  
In this version we need to sort a list and then iterate over it: the cost of this is not too high because this list contains 24 elements at most.
- (4) FlatMap to obtain an RDD like  $((Date, ID), variation)$
- (5) Remaining steps are the same as before from step 8 onwards

Thus the two variants differ only in the way they compute the last price before of an hour.

## 5.2 Query 2 with Spark SQL

As for this query using DataFrame we did the following steps:

- (1) Find the last trade info for each hour  $hh:00$  and then, summing 1 to the hour, find the first price for  $hh:00$
- (2) Find price couples looking for, for each hour  $hh:00$ , to the most recent previous trade
- (3) Compute for each hour the price change, obtaining a schema of  $(Date, ID, Variation)$
- (4) Compute statistics grouping by  $Date$  and  $ID$
- (5) Look for best and worst stocks for each date using nested queries with windows and *row\_number()*
- (6) Union between best and worst stocks data

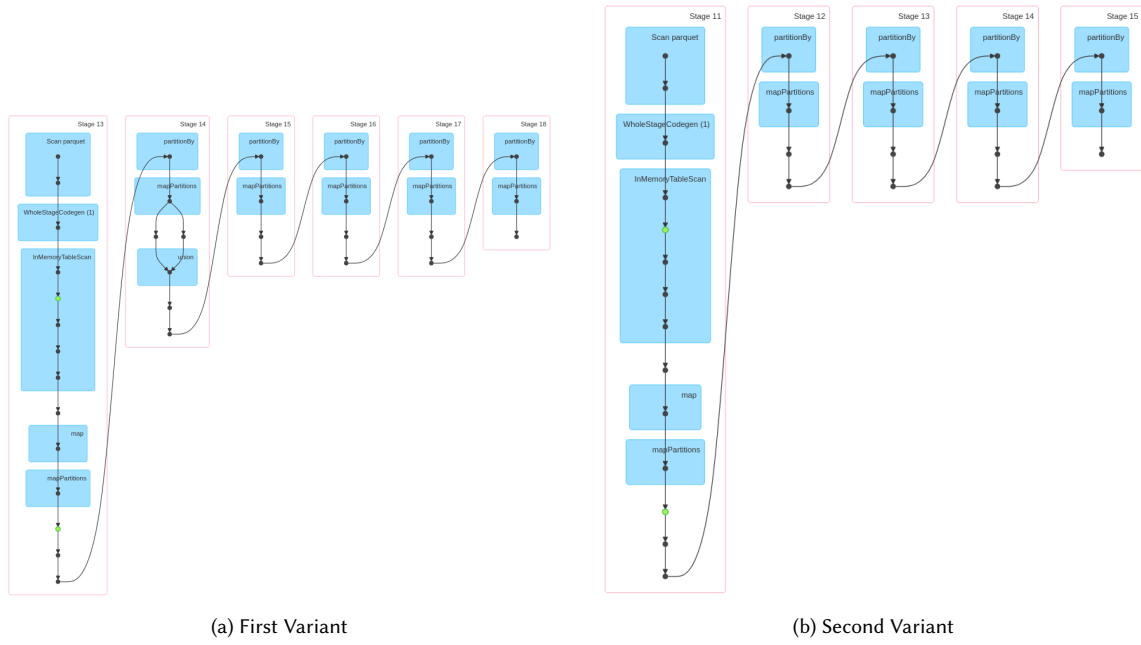


Fig. 3. DAGs for Query 2 with RDD

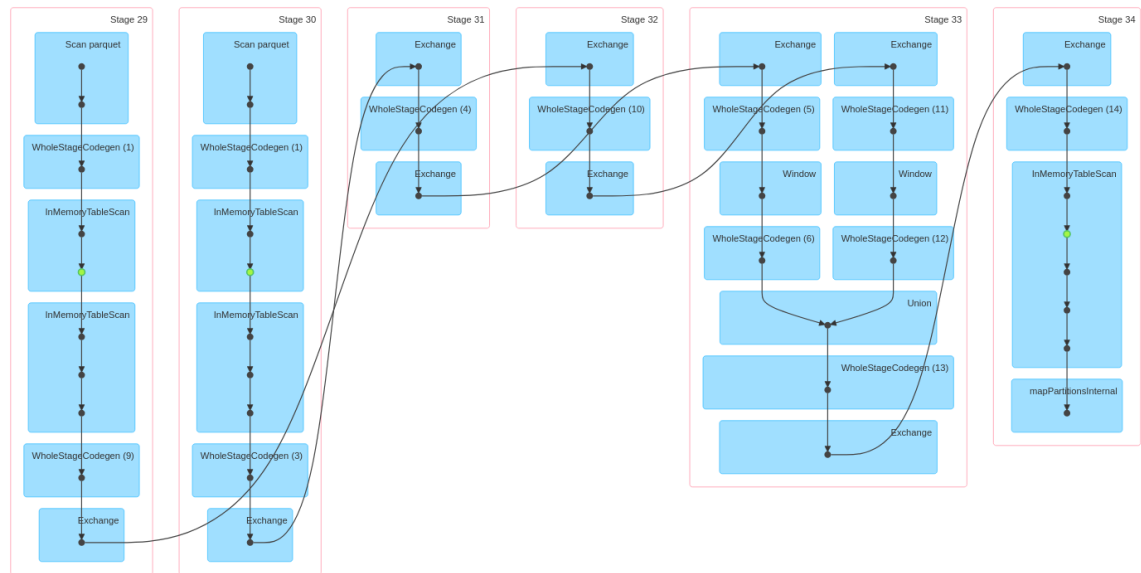


Fig. 4. DAGs for Query 2 with DataFrame

## 6 QUERY 3

**Query** For each day, calculate the 25th, 50th, 75th percentile of the change in the price of sale of stocks traded on individual markets. The statistic should be calculated by considering all and only the stocks belonging to each specific market: Paris (FR), Amsterdam (NL) and Frankfurt/Xetra (ETR). In the output also indicate the total number of events used to calculate the statistics.

### 6.1 Query 3 with RDD

With regard to this query using RDD we did the following steps:

- (1) Map to a key-value RDD like  $((Date, ID), (Time, Last, Time, Last))$
- (2) ReduceByKey in order to find, for each couple  $(Date, ID)$  the first and last price of the day obtaining an RDD like  $(Date, ID) (firstTime, firstPrice, lastTime, lastPrice)$
- (3) Filter for those elements with  $firstTime \neq lastTime$  in order not to consider those stocks with only one trade during the day
- (4) Map to find variation doing the difference between found prices
- (5) Map to extract from every  $ID$  the market of belonging obtaining an RDD with key  $(Date, Market)$  obtaining an RDD like  $(Date, Market) (Variation)$
- (6) GroupByKey to find for each  $(Date, Market)$  the list of variations
- (7) Map to sort this list and then find percentiles.

We computed percentiles using a custom function named *computePercentile* which takes as input a list and a float for the percentile value. This function multiplies list length with percentile value and then if the product is an integer, it computes the percentile as the average of two adjacent elements of the list, otherwise it returns the element whose index is the rounded product value.

### 6.2 Query 3 with Spark SQL

As for this query using DataFrame we followed the same logic as the RDD version:

- (1) Find, using a nested query, for each couple  $Date, ID$ , the first and last trade time excluding those with same min and max time as they had no variation during that date.
- (2) Find using a double *Join* the first and last price of each Stock
- (3) Find price change over the day
- (4) Extract Market from  $ID$
- (5) Find percentiles grouping by  $Date$  and  $Market$  and then aggregating with *approx\_percentile*

## 7 EMPIRICAL RESULTS

All execution times have been taken without considering either pre process times to convert the read file and post processing time to prepare results for saving on local file system, Redis and HDFS. Times have been taken using Python *time* library, taking time before and after a *collect()* on the result. All experiments have been done using Docker Compose on a single node.

We executed ten runs for each query; after the execution of any of the three queries we closed both *SparkContext* and *SparkSession* in order to reinitialize execution context and not have influences between executions. We considered Spark Workers configured as follows:

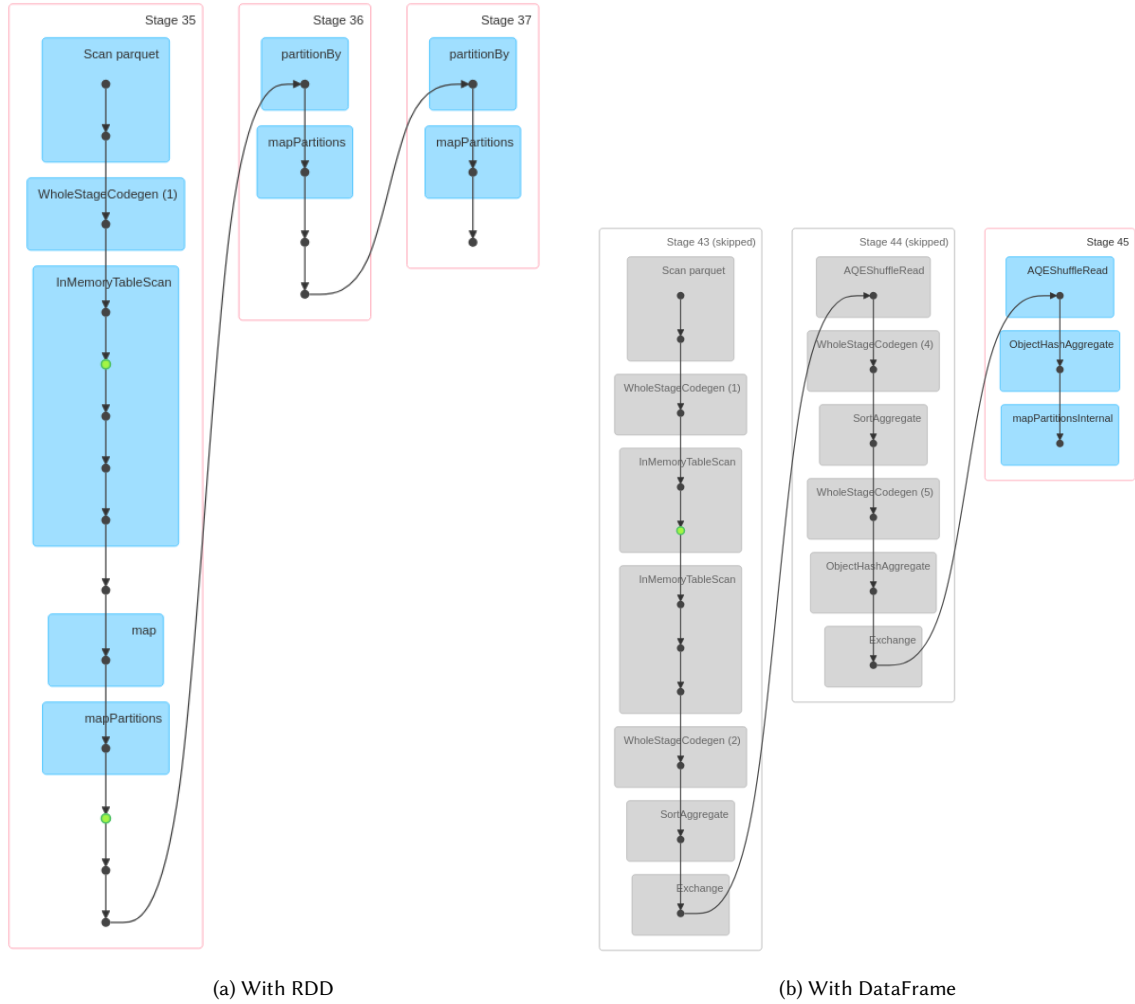


Fig. 5. DAG for Query 3

- SPARK\_WORKER\_MEMORY=1G
- SPARK\_WORKER\_CORES=1

As we can see from the charts and from tables of values, a query executed using RDD is nearly always faster than the same query executed using SQL, except for Query 1, where SQL wins by a little difference. This higher RDD speed is due to the fact that RDD is a lower abstraction in Spark and, as a consequence, does not suffer the overhead as the DataFrame does.

As for queries implemented using RDD, we can see how increasing the number of Spark workers does not affect significantly execution time; even in Query 1, when we have the highest variation, considering the scale it is not that much.



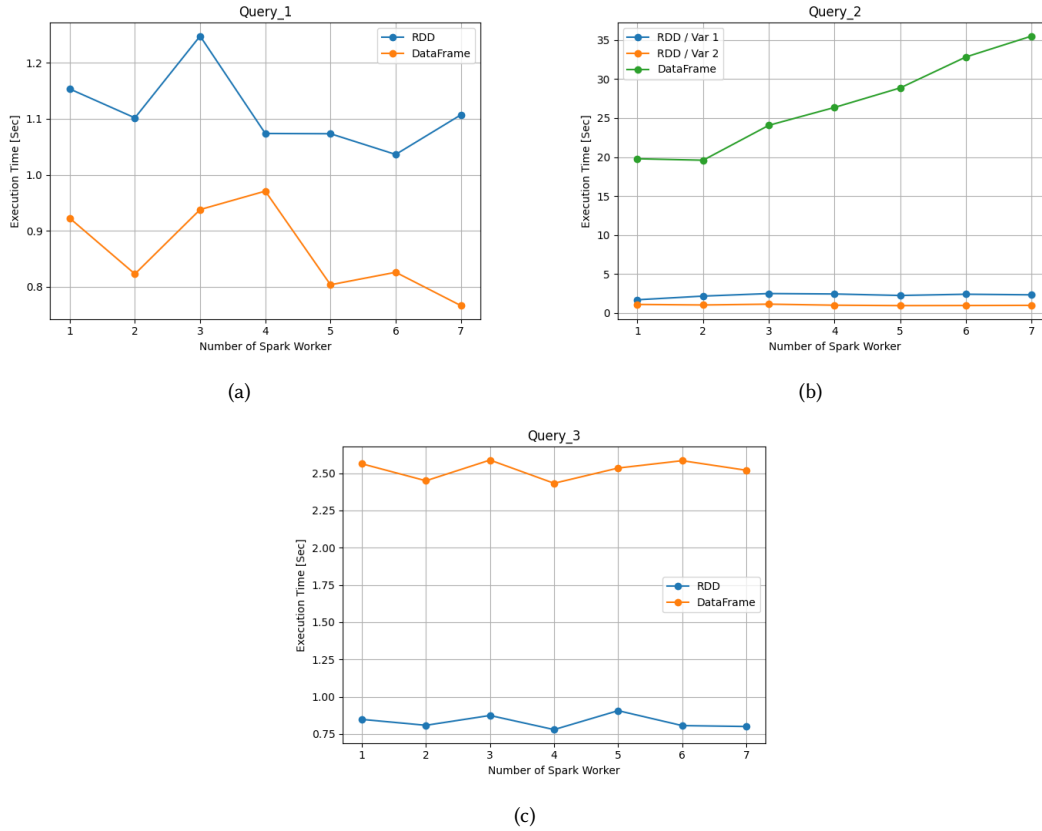


Fig. 6. Execution Times Charts

As for Query 2 implemented using RDD, we can see how in the first version there is a slow increase in execution time, while in the latter there is a more or less constant trend: this is because the first version uses *join* which with one node is executed locally, while with multiple nodes requires shuffling between workers. As a consequence we can say that the small gap between the two versions is the time for *Join* execution, as the two versions differ only for this operation and few others until obtaining the variations list.

Even for queries implemented using Spark SQL we can see how increasing the number of Spark workers does not affect significantly execution time except for Query 2, when we have an high number of *Joins* and consequently an high number of shuffles between workers.

Speaking generally, we can see that since the filtered and cleaned dataset is quite small, there is no significant improvement by increasing the number of workers.

Finally we can say that these charts show only a partial vision of time trend to the increment of workers: indeed, due to the fact that the execution is on a single node, the network latency is near zero and it is not considered in these charts or it is poorly represented.

<b>DataStructure</b>	<b>Variant</b>	<b>WorkersNum</b>	<b>Avg</b>
RDD	1	1	1.6977777481079102
RDD	1	2	2.169368124008179
RDD	1	3	2.4952018737792967
RDD	1	4	2.4392569780349733
RDD	1	5	2.250243401527405
RDD	1	6	2.4078279972076415
RDD	1	7	2.332538366317749
RDD	2	1	1.1006942987442017
RDD	2	2	1.0299310445785523
RDD	2	3	1.1350558996200562
RDD	2	4	1.0022500038146973
RDD	2	5	0.9502127647399903
RDD	2	6	0.9629277944564819
RDD	2	7	0.9825979948043824
DataFrame	1	1	19.781137919425966
DataFrame	1	2	19.59350690841675
DataFrame	1	3	24.072310662269594
DataFrame	1	4	26.363581919670104
DataFrame	1	5	28.871790885925293
DataFrame	1	6	32.840597248077394
DataFrame	1	7	35.51872515678406

Table 2. Query 2 Averages

<b>DataStructure</b>	<b>Variant</b>	<b>WorkersNum</b>	<b>Avg</b>
RDD	1	1	1.1534550189971924
RDD	1	2	1.101752471923828
RDD	1	3	1.2479825735092163
RDD	1	4	1.0738351583480834
RDD	1	5	1.0734047651290894
RDD	1	6	1.0363578796386719
RDD	1	7	1.107157588005066
DataFrame	1	1	0.9225053548812866
DataFrame	1	2	0.822759461402893
DataFrame	1	3	0.9378469705581665
DataFrame	1	4	0.9708062410354614
DataFrame	1	5	0.8034405946731568
DataFrame	1	6	0.8257670879364014
DataFrame	1	7	0.7662054300308228

Table 3. Query 1 Averages

<b>DataStructure</b>	<b>Variant</b>	<b>WorkersNum</b>	<b>Avg</b>
RDD	1	1	0.8482024431228637
RDD	1	2	0.8086688280105591
RDD	1	3	0.8745376586914062
RDD	1	4	0.7797282934188843
RDD	1	5	0.9067367792129517
RDD	1	6	0.806633186340332
RDD	1	7	0.800520944595337
DataFrame	1	1	2.562646722793579
DataFrame	1	2	2.448776030540466
DataFrame	1	3	2.5877469778060913
DataFrame	1	4	2.4319807291030884
DataFrame	1	5	2.5341213941574097
DataFrame	1	6	2.583537793159485
DataFrame	1	7	2.5193403482437136

Table 4. Query 3 Averages